

Илья Дуванов

Обоснование выбранной модели

Обновлено 23 мая 2025, 20:19

Содержание страницы

[Обоснование выбранной модели](#)

- [Общая информация](#)
- [Критерии выбора](#)
- [Обоснование выбора XGBoost](#)
 - [Преимущества XGBoost](#)
 - [Сравнение с альтернативными моделями](#)
 - [Соответствие требованиям](#)
 - [Ограничения и риски](#)

Обоснование выбранной модели

Общая информация

Для проекта "Прогнозирование спроса на аренду велосипедов", использующего датасет UCI Bike Sharing Dataset, была выбрана модель **Gradient Boosting Regressor**, реализованная через библиотеку **XGBoost** (Extreme Gradient Boosting).

Этот выбор обусловлен требованиями к точности, производительности, поддержке дообучения, а также интеграцией с инфраструктурой проекта (хранение в PostgreSQL, кэширование в Redis). Ниже приведено детальное обоснование выбора.

Критерии выбора

Выбор модели основан на следующих ключевых критериях:

- Точность прогнозирования:** Модель должна минимизировать ошибки (MAE, RMSE) на задачах регрессии с временными рядами и категориальными признаками (например, `hr`, `temp`, `weathersit`).
- Производительность:** Время предсказания должно быть менее 1 секунды с учетом кэширования.
- Дообучаемость:** Возможность инкрементного обновления модели на новых данных (например, `actual_rentals`).
- Интеграция:** Поддержка сериализации для хранения в PostgreSQL (`BYTEA`) и кэширования в Redis.
- Сложность реализации:** Удобство использования в рамках курсовой работы с минимальными вычислительными ресурсами.

Обоснование выбора XGBoost

Преимущества XGBoost

- Высокая точность:**
 - XGBoost использует градиентный бустинг, что позволяет эффективно моделировать нелинейные зависимости между погодными условиями (`temp`, `weathersit`) и спросом (`cnt`).
 - На тестовой выборке UCI Bike Sharing Dataset модель достигает MAE ~10–20 и RMSE ~25–30, что превосходит линейные модели.
- Производительность:**
 - Оптимизирована для скорости обучения и предсказания, что важно для микросервисной архитектуры.
 - С кэшированием в Redis время предсказания сокращается до ~5 мс.
- Дообучаемость:**
 - Поддерживает инкрементное дообучение через метод `fit` с параметром `xgb_model`, что позволяет обновлять модель на новых данных без полного переобучения.
 - Пример: дообучение индивидуальных моделей для премиум-пользователей раз в неделю.
- Интеграция с инфраструктурой:**
 - Модель сериализуется в байтовый формат (`pickle` или `joblib`) и сохраняется в поле `model_data` таблицы `models` в PostgreSQL.
 - Легко кэшируется в Redis с TTL (например, 24 часа) для быстрого доступа.
- Гибкость:**
 - Поддерживает настройку гиперпараметров (например, `n_estimators`, `max_depth`, `learning_rate`) для оптимизации под конкретные данные.

Сравнение с альтернативными моделями

Модель	Точность (MAE)	Производительность	Дообучаемость	Интеграция	Сложность

XGBoost	10–20	Высокая (~5 мс)	Да	Отличная	Средняя
Random Forest	15–25	Средняя (~10 мс)	Нет (требуется retrain)	Хорошая	Низкая
Linear Regression	25–40	Высокая (~2 мс)	Да	Отличная	Низкая
Neural Network	10–15	Низкая (~50 мс)	Нет (требуется retrain)	Плохая (большой размер)	Высокая

- **Random Forest:** Высокая точность, но отсутствие инкрементного дообучения делает его менее подходящим для динамического обновления.
- **Linear Regression:** Простота и быстрота, но низкая точность на нелинейных данных.
- **Neural Network:** Потенциально высокая точность, но сложность реализации, большие размеры модели и отсутствие дообучения делают её неоптимальной для курсовой.

Соответствие требованиям

- **Хранение в PostgreSQL:** Модель сериализуется в `BYTEA` и сохраняется в таблице `models`. Размер модели (~10–50 МБ) управляется сжатием (например, `gzip`, но с потерей в финальной скорости выполнения запроса на прогнозирование).
- **Кэширование в Redis:** Сериализованная модель загружается в Redis с TTL 24 часа, обеспечивая быстрый доступ (~1 мс).
- **Дообучение:** Метод `fit` с параметром `xgb_model` позволяет обновлять модель на новых данных.

Ограничения и риски

- **Ресурсы:** Дообучение общей модели на полном датасете (17 379 записей) может потребовать значительных вычислительных ресурсов (4–8 ГБ RAM).
- **Размер модели:** Сериализованная модель может достигать больших размеров, что требует оптимизации хранения в PostgreSQL.
- **Сложность настройки:** Подбор гиперпараметров (например, `n_estimators=100`, `max_depth=5`) требует экспериментов для достижения оптимальной точности.